

STORAGE NETWORK CONTROL METHOD

Abstract of Prior Art

Publication number: JP10069357

Publication date: 1998-03-10

Inventor: TAKAMOTO YOSHIFUMI; KANAI HIROKI

Applicant: HITACHI LTD

Classification:

- international: G06F3/06; G06F12/00; G06F13/00; G06F3/06;
G06F12/00; G06F13/00; (IPC1-7): G06F3/06; G06F3/06

- European:

Application number: JP19960226401 19960828

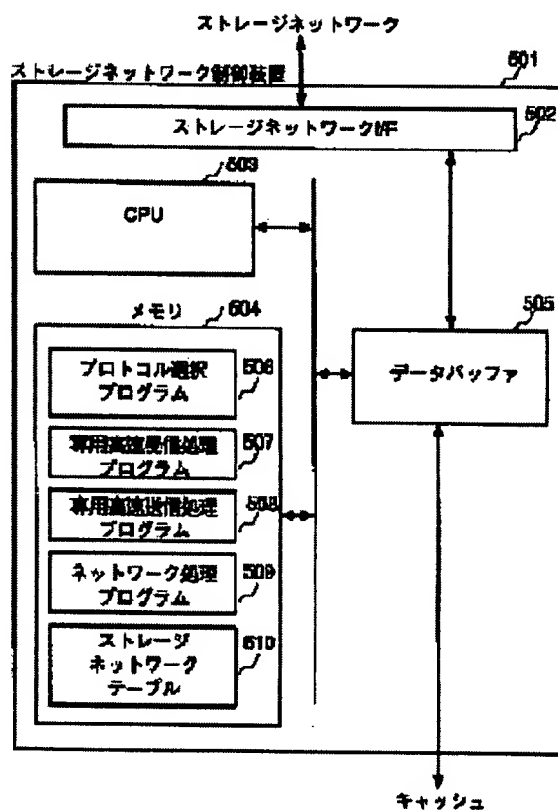
Priority number(s): JP19960226401 19960828

Report a data error here

Abstract of JP10069357

PROBLEM TO BE SOLVED: To perform communication at a high speed between a disk array controllers while keeping a general purpose interface with a host by providing a mechanism for switching a communication protocol between the host and the disk array controller and the communication protocol between the disk array controllers.

SOLUTION: Information relating to devices connected on a storage network is stored in a storage network table 510. A protocol selection program 506 selects dedicated high-speed reception/transmission processing programs 507 and 508 or a network processing program 509 from the contents of the storage network table 510 and switches the protocol to be called corresponding to the result. Thus, the network processing program 509 is executed as before when a communication destination is the host and the dedicated high-speed reception/transmission processing programs 507 and 508 are executed when it is the disk array controller.



Data supplied from the esp@cenet database - Worldwide

Prior Art 1.

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11)特許出願公開番号

特開平10-69357

(43)公開日 平成10年(1998)3月10日

(51) Int.Cl. ^o	識別記号	序内整理番号	F I	技術表示箇所
G 0 6 F 3/06	3 0 1		G 0 6 F 3/06	3 0 1 P 3 0 1 X
	5 4 0			5 4 0

審査請求 未請求 請求項の数3 O.L (全 10 頁)

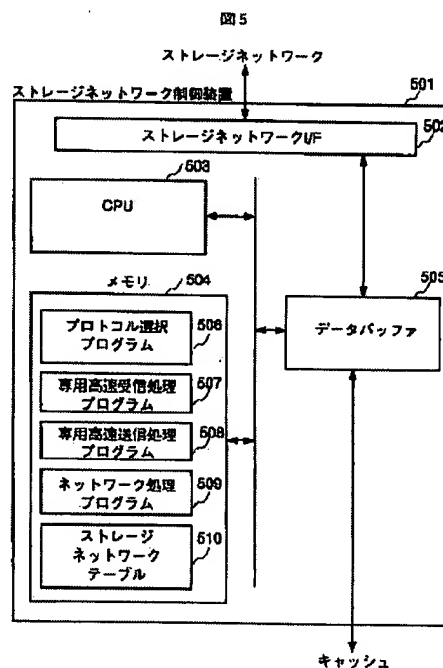
(21)出願番号	特願平8-226401	(71)出願人	000005108 株式会社日立製作所 東京都千代田区神田駿河台四丁目6番地
(22)出願日	平成8年(1996)8月28日	(72)発明者	高本 良史 東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所中央研究所内
		(72)発明者	金井 宏樹 東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所中央研究所内
		(74)代理人	弁理士 小川 勝男

(54)【発明の名称】 ストレージネットワーク制御方法

(57) 【要約】

【課題】ディスクアレイ制御装置とホストプロセッサの通信プロトコルと、ディスクアレイ制御装置間の通信プロトコルを切り替える。

【解決手段】ディスクアレイ制御装置内にストレージネットワーク制御装置501を設け、この中に設けたストレージネットワークに接続された装置がホストプロセッサかディスクアレイ制御装置かを示すストレージネットワークテーブル510をサーチし、送信先に応じて通信プロトコルを切り替える。



【特許請求の範囲】

【請求項1】複数のディスク装置を並列に動作させるとともに冗長データを格納する複数のディスクアレイ制御装置と、単一あるいは複数のホストプロセッサを接続するストレージネットワークを有するシステムにおいて、前記ストレージネットワークに接続された装置が、前記ホストプロセッサかディスクアレイ制御装置か、を示すストレージネットワークテーブルを前記ディスクアレイ制御装置内に設け、前記ディスクアレイ制御装置はデータの送信時に前記ストレージネットワークテーブルをサーチし、送信先が前記ホストプロセッサか前記ディスクアレイ制御装置かを判定し、前記ディスクアレイ制御装置であれば簡略化した通信方法を使用して通信を行うことを特徴とするストレージネットワーク制御方法。

【請求項2】請求項1において、前記ディスクアレイ制御装置および前記ホストプロセッサはデータ送信時に、送信元が前記ディスクアレイ制御装置か、前記ホストプロセッサかを示す識別子を送信データに付加し、前記ディスクアレイ制御装置は受信時に前記識別子を参照し、送信元が前記ディスクアレイ制御装置であれば前記簡略化した通信方法を使用するストレージネットワーク制御方法。

【請求項3】請求項1において、前記ストレージネットワークに接続された装置は、起動時に自装置識別子を送信し、前記自装置識別子を受信した装置は前記ストレージネットワークテーブルに追加するストレージネットワーク制御方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、ディスク制御装置に係り、ネットワークで相互に接続されたディスク制御装置のデータ通信制御方法に関する。

【0002】

【従来の技術】一般にコンピュータシステムはプロセッサと二次記憶装置から構成されている。主として使用される二次記憶装置は磁気ディスク装置である。現在、ディスク記憶装置の容量の伸び率は極めて高いが、メカニカルな動作を伴う磁気ディスク装置の性能はプロセッサ性能の伸び率ほど高くない。その課題を解決する方法として、以下に示すディスクアレイが提案された。

【0003】ディスクアレイは、磁気ディスクを並列に制御することで性能を向上させる方式である。D. Patterson, G. Gibson, and R. H. Kartzらによる“A Case for Redundant Arrays of Inexpensive Disks (RAID), in ACM SIGMOD Conference, Chicago, IL”, pp. 109-116 (June 1988) (以下、第1の文献と呼ぶ)では、複数のディスクドライブにデータを分散して配置することでディスク内に格納されたデータへのアクセス時間を短縮し、かつパリティあるいはECCと呼ばれる冗長データを格納することで信頼性も高めるRAIDというディス

クアレイの構成技術が紹介されている。つまり、アレイディスクでは、複数のディスクドライブに対して並列に入出力を行うので、データの読み出しあるいは書き込みは高速となり、また、ディスクドライブに障害が発生したときでもパリティと障害ディスクドライブ以外のデータから、障害ディスクドライブのデータを回復することができる。

【0004】より具体的には、上記第1の文献では、データの格納方法によりRAIDレベルを複数に分類している。そのうち、製品で多く使用されるRAIDレベルは、RAID1, RAID3, RAID5である。

【0005】RAID1はミラーリングであり、2台のディスクに同一データを格納することでディスクの障害に対する信頼性を高めている。読み出し時には、2台のディスクの内どちらか早いほうのディスクから読むことで単一ディスクに比べ高速である。RAID3は、単一データをビットあるいはバイト毎にストライピング(分割)し、複数のディスクに並列に読み書きすることで、データ転送性能を向上可能である。RAID5は入出力ブロックを単位として複数のディスクにストライピングを行うレベルである。

【0006】RAID3では単一入出力要求を小さな単位でストライピングするのに対し、RAID5はそれよりも大きな単位でストライピングを行う。RAID3は単一ユーザの大規模データ入出力を高速化することが主な目的であるが、RAID5は多数ユーザの小規模データ入出力を並列に実行することが主な目的である。従って、RAID3は特に大規模なデータ入出力が要求されるマルチメディアや科学技術計算に用いられる。RAID5は多数ユーザのサービスを行うオンライン・トランザクションなどでのデータベース処理に用いられる。

【0007】しかし、単一のディスクアレイでは性能や容量の増加に限界が生じる。そのために、ディスクアレイをさらに複数設け、並列度をより増やす方法がある。大容量で低コストのディスクアレイを構築する場合、単一のディスクアレイコントローラに多数のディスク装置を接続する構成が考えられるが、ディスクアレイ制御装置内の制御プロセッサの性能の限界のために容量は増加するが性能は伸びなくなってしまうという問題が発生する。

【0008】制御プロセッサはホストからの入出力要求を解釈実行するために設けられる。制御しなければならぬディスク装置の台数に比例した処理性能が要求されることから、現状では数台から数十台規模のディスクアレイがほとんどである。

【0009】この問題を解決する方法として、バスあるいは高速なネットワークで制御装置を接続することで拡張性を高める方法が考えられる。多数のディスクアレイを接続する方法は特開平7-44322号公報に述べられている。これは制御装置を階層化することで、各制御装置が

制御しなければならない台数を削減することが目的である。例えば、一つの制御装置で5台のディスク装置を制御できるとする。この制御装置を単一のバスあるいはネットワークに接続し、ホストプロセッサからはいずれの制御装置へもアクセスできるようにする。こうすることで全体で多数台のディスク装置を制御することが可能になる。

【0010】

【発明が解決しようとする課題】ディスクアレイはホストプロセッサと接続され、ホストプロセッサからディスクアレイ制御装置に対して読み込みあるいは書き込みコマンドが発行されると、それに対応したデータの読み込みあるいは書き込みを行い、その実行結果をホストプロセッサに返送する。このような使用環境では、通信の形態は単一のホストプロセッサと複数のディスクアレイ制御装置との間で行われる。前述の特開平7-44322号公報がその例である。

【0011】しかし、ディスクアレイ制御装置間で自動的にバックアップの取得を行うような処理の場合、複数のディスクアレイ制御装置が接続されたバスあるいはネットワークを介して直接データを転送することになる。この場合問題となるのは通信方法（通信プロトコル）である。ディスク制御装置がホストプロセッサと通信を行う場合には、ホストプロセッサのインタフェースに合わせた通信プロトコルでデータの送受信を行う。しかし、ホストプロセッサのインタフェースは汎用である分、性能は低い。これは、ホストプロセッサが、ソフトウェアでデータのエラーチェックを行っていることが大きな要因である。

【0012】ホストのインタフェースに合わせた通信プロトコルしか有さないディスクアレイ制御装置を使用して、前述のディスクアレイ制御装置間のバックアップを行った場合も同様の理由でバックアップ性能が低下することになる。これは、単一の通信プロトコルしか有していないことと、通信プロトコルを切り替えるための手段を有していないことが理由である。特開平7-44322号公報では、複数のディスクアレイ制御装置を接続する形態が述べられているが、ディスクアレイ制御装置間で通信を行うことについては述べられていない。

【0013】

【課題を解決するための手段】上記問題を解決するために、複数のディスク装置を並列に動作させるとともに冗長データを格納する複数のディスクアレイ制御装置と、単一あるいは複数のホストプロセッサを接続するストレージネットワークを有するシステムにおいて、前記ストレージネットワークに接続された装置が、前記ホストプロセッサかディスクアレイ制御装置か、を示すストレージネットワークテーブルを前記ディスクアレイ制御装置内に設け、前記ディスクアレイ制御装置はデータの送受信時に前記ストレージネットワークテーブルをサーチし、

送信先が前記ホストプロセッサか前記ディスクアレイ制御装置かを判定し、前記ディスクアレイ制御装置であれば簡略化した通信方法を使用して通信を行う。これにより、ホストプロセッサとの通信時は汎用インタフェースによる通信プロトコルを使用し、ディスクアレイ制御装置間で通信を行う場合には高速な通信プロトコルを使用することができる。

【0014】

【発明の実施の形態】以下、本発明のストレージネットワーク制御方法を図面に示し実施例を参照してさらに詳細に説明する。

【0015】図1は本発明によるストレージネットワーク制御方法を適用するディスク装置のブロック図である。101、102はホストであり、106、107はディスクアレイ制御装置であり、105はストレージネットワークである。ホストプロセッサ101、102はストレージネットワーク105を介してディスクアレイ制御装置106、107に相互接続されている。130、131はディスク装置であり、それぞれ複数のディスクドライブで構成されている。

【0016】ホストプロセッサ101、102から発行されたディスクアレイ入出力要求は、ストレージネットワーク105を介してディスクアレイ制御装置106、107に転送され、またディスクアレイ入出力要求に伴うデータもストレージネットワーク105を介して転送される。従来から存在する主にコンピュータ間のデータ通信を主な目的とするネットワーク132は、ストレージネットワーク105とは独立に接続されている。ホストプロセッサ101、102内には、ストレージネットワーク105の制御を行うストレージネットワークインタフェース103、104が設けられる。

【0017】ディスクアレイ制御装置106、107内には、ネットワークストレージ105を制御するインタフェース制御部106、107、ディスクアレイ制御装置106、107を制御するCPU110、111とマイクロプログラム118、119、ディスク装置130、131のデータを一時的に格納するキャッシュ112、113、ディスクコマンドを分配するスイッチ116、117、ディスク装置130、131との接続を制御するインタフェース制御120～124、125～129、パリティ制御114、115から構成されている。

【0018】ディスクアレイ制御装置106、107の主な機能は、複数のディスクドライブを使用し、性能と信頼性を向上させるディスクアレイ制御を行うことである。ディスクアレイは、データやパリティの格納方法により、いくつかに分類されているが、本実施例では、それらのいずれも使用可能であり、また混在も可能である。それらの主な種類には前述の通り、RAID1、RAID3、RAID5などがある。

【0019】RAID1、3、5ではいずれの場合も、

データ格納時に冗長データも同時にディスクに書き込むことで信頼性も向上している。具体的には、RAID1の場合は、同一データを複数のディスクに書き込むことで冗長性を持たせている。またRAID3、5では、データを複数のディスクに分割するが、この時、旧データと旧パリティと新データとの排他的論理和を新パリティとしてディスクに書き込む。こうすることで、いずれかのディスクが障害を起こしても、RAID1、3、5いずれの場合も障害を起こしたディスクのデータを回復することが可能である。

【0020】例えば、RAID3において以下のようにデータが格納されているとする。

ディスク1 データ = 1 0 1 0 1 0 1 0

ディスク2 データ = 1 1 1 1 0 0 0 0

ディスク3 パリティ = 0 1 0 1 1 0 1 0

この状態でディスク2が障害を起こし読み書きができなくなった場合、ディスク1とディスク3の排他的論理和を演算する。

ディスク1 データ = 1 0 1 0 1 0 1 0

ディスク3 パリティ = 0 1 0 1 1 0 1 0

論理和
 = 1 1 1 1 0 0 0 0 (ディスク2のデータ)

このように、ディスク2のデータを再現することが可能である。図1のパリティ制御114、115はパリティの生成あるいはデータの回復を行う機構である。

【0021】本発明は、ストレージネットワークの通信方法に関するが、図2を用いてストレージネットワークの特徴の一例を示す。ホストプロセッサ201、202はストレージネットワーク204を介してディスクアレイ制御装置205、206に相互接続されている。207、208はディスク装置であり、それぞれ複数のディスクドライブで構成されている。図2はディスクアレイ制御装置205からディスクアレイ制御装置206へデータのバックアップを行う動作を示している。ストレージネットワーク204を設けることで、ディスクアレイ制御装置205、206間で直接バックアップを取得することが可能となる(209)。

【0022】従来のバックアップ(図示せず)は、ホストプロセッサがバックアップ元のディスクアレイ制御装置からバックアップデータの一部分を読み込み、そのデータをバックアップ先のディスクアレイ制御装置へネットワーク等を介して転送する操作を繰り返し行うことで処理していた。この従来の方法では、バックアップデータのすべてをホストプロセッサが処理しなくてはならず、バックアップ量が膨大になるとホストプロセッサはバックアップのために相当の能力を使用しなければならなくなる。しかし、ストレージネットワーク204を設けることで、ホストプロセッサはバックアップ要求203を一回発行するだけで、バックアップはディスクアレイ制御装置205、206間で自動的に取得することが

でき、またホストプロセッサ201、202へ処理負荷を与えることもなくなる。

【0023】図3は本発明の概要を示した図である。ホストプロセッサ301、302はストレージネットワーク303を介してディスクアレイ制御装置306、307に相互接続されている。310、311はディスク装置であり、それぞれ複数のディスクドライブで構成されている。ストレージネットワーク303は、ディスクアレイ制御装置306、307とホストプロセッサ310、302間で通信すること以外に、ディスクアレイ制御装置306、307間でも通信可能である。

【0024】本発明の特徴は、ホストプロセッサ301、302とディスクアレイ制御装置306、307間で通信する時と、ディスクアレイ制御装置306、307間で通信する時の通信プロトコルを自動的に切り替えることである。ホストプロセッサ301、302とディスクアレイ制御装置306、307間で通信する時にはネットワークプロトコルを使用し(304)、ディスクアレイ制御装置306、307間で通信する時は専用高速プロトコルを使用する(305)。

【0025】ここで、プロトコルとは相互に情報を送受するための取り決められた通信方法である。通信におけるプロトコルは階層分けされている。一般には、物理層と論理層に分けられる。物理層は、データを運ぶためのケーブルの電気的特性等が決められている。また論理層では、ユーザデータを転送するためのソフトウェア的な通信手順が決められている。階層分けすることで、論理層が同一であれば物理層だけ異なるメディアを用いても通信可能となる。

【0026】図3で示したネットワークプロトコルと専用高速プロトコルは論理層に対応している。ホストプロセッサ301、302を使用するユーザにストレージネットワークを意識させる必要がないようにするために、ユーザは従来のネットワーク132と同一の論理層のプロトコルを用いてストレージネットワークを使用することができる。これは、後で詳細に説明する。ユーザからは従来のネットワーク132もストレージネットワーク303もまったく同じ操作で使用可能である。これは、前述の論理層のプロトコルを従来のネットワーク132と同一のプロトコルに合わせたことの利点である。

【0027】しかし、ネットワークプロトコル304は相互接続性では優れているが、性能面は高くない。これは、汎用性を持たせることと相反することであり、ハードウェア依存を少なくして通信できるようにするためにはハードウェアの機能を削減し、代わりにソフトウェアの機能を高くすることで実現されているためである。具体的には、ソフトウェアで転送時のエラーチェック機能を行っていることが性能低下の大きな原因になっている。これに対し、専用高速プロトコル305はエラーチェック機能を省略することで高速化が可能となってい

る。従来のネットワークとの互換はなくなるが、ユーザに見えない部分での通信であるためユーザへの影響はない。

【0028】図4はホストプロセッサ401の構造を示している。ここでは、従来のネットワーク132とストレージネットワーク105がどのように共存しているかを示している。ホストプロセッサ401は、演算装置402、メモリ403、ネットワークインタフェース408、ストレージネットワークインタフェース409から構成されている。

【0029】演算装置402はメモリ403から命令を読み取り解釈／実行を行う。メモリ403内には、ユーザプログラム404、ネットワーク処理プログラム405、ネットワークドライバ406、ストレージネットワークドライバ407が格納されている。ユーザプログラム404は、ネットワーク処理プログラムに対して入出力要求を発行すると、要求内容によってネットワークドライバ406あるいはストレージネットワークドライバ407に処理を切り替えて要求を発行する。それを受けたネットワークドライバ406あるいはストレージネットワークドライバ407は、異なるハードウェアであるネットワークインタフェース408とストレージネットワークインタフェース407をそれぞれ制御し、接続相手と通信を行う。ネットワーク処理プログラムは図3で述べた、論理層のプロトコルを処理する。ユーザは、通信相手がネットワークかストレージネットワークかを意識することなくプログラムを開発し使用することが可能となる。

【0030】図5は本発明の特徴であるストレージネットワーク制御を行う制御装置の構成を示している。ストレージネットワーク制御装置501はストレージネットワークインタフェース502、CPU503、メモリ504、データバッファ505から構成されている。ストレージネットワークインタフェース502はストレージネットワーク105に接続するための電氣的な制御を行う。CPU503はストレージネットワーク制御装置501の全体を制御する。具体的には、メモリ504に格納されているプログラムを随時読み出し、解釈／実行を行う。メモリ504内には、プロトコル選択プログラム506、専用高速受信処理プログラム507、専用高速送信プログラム508、ネットワーク処理プログラム509、ストレージネットワークテーブル510が格納されている。

【0031】プロトコル選択プログラム506は、専用高速受信／送信処理プログラム507、508かネットワーク処理プログラム509かのいずれかを選択し、呼び出す処理を行う。この処理については後で詳細に説明する。専用高速送信／受信処理プログラム507、508はネットワーク処理プログラム509と比較し汎用性はないが、処理が高速であることが特徴である。この処理

については、後で詳細に説明する。

【0032】ストレージネットワークテーブル510はストレージネットワーク上に接続された装置に関する情報が格納されている。プロトコル選択プログラム506は、ストレージネットワークテーブル510の内容から、どの処理プログラムを選択するかを決定する。データバッファは、キャッシュあるいはストレージネットワークから転送されたデータを一時的に格納するために使用される。また、専用高速送信／受信処理プログラム507、508やネットワーク処理プログラム509がデータにヘッダを付加するためにも使用される。

【0033】図6、図7はストレージネットワークに接続された装置の例と、その場合のストレージネットワークテーブル510の例を示している。

【0034】図6では、ストレージネットワークに2台のホスト601、602と2台のディスクアレイ制御装置604、605が接続されている。2台のホスト601、602と2台のディスクアレイ制御装置604、605はストレージネットワーク603で相互通信が可能な構成で接続されている。2台のホスト601、602は、2台のディスクアレイ制御装置604、605のいずれからでもデータ入出力が可能である。また、ディスクアレイ制御装置604、605間でもデータ送受信が可能である。ホスト601、602は、それぞれネットワークストレージ識別子1、2が定義されている。また、ディスクアレイ制御装置604、605にも同様にネットワークストレージ識別子3、4が定義されている。

【0035】図7はこの構成例における、ストレージネットワークテーブル510内の構造を示している。テーブルのカラム701はストレージネットワーク識別子が格納されている。テーブルのカラム702はストレージネットワーク識別子に対応した装置の種類が格納されている。例えば、ホストであれば0が、ディスクアレイ制御装置であれば1が格納されている。プロトコル選択プログラム506は、このテーブルの内容から、要求されたコマンドやデータをどのような装置に転送するのかサーチし、その結果に応じて呼び出すべきプロトコルを切り替える。具体的には、通信先がホストであれば従来どおりネットワーク処理プログラムを実行し、ディスクアレイ制御装置であれば専用高速送信／受信処理プログラムを実行する。

【0036】ストレージネットワークテーブル510の生成は、システム起動時にストレージネットワークに接続された装置間で情報の受け渡しを行うことで実現可能である。具体的には、自システム起動時にストレージネットワークに接続された全ての装置に対し自装置の識別子を発行する。識別子の転送は識別子転送用のコマンドを使用する。このコマンドを受けた装置は、自装置内のストレージネットワークテーブルに送信された識別子を追加する。

【0037】図8はプロトコル選択プログラムの処理フローを示している。

【0038】ステップ801では要求が送信要求か受信要求かを判断する。送信要求であればステップ802に移り、受信要求であればステップ807に移る。ステップ802では送信先識別子を取得する。これはキャッシュ108、109からデータバッファ505内に転送されたデータにデータと共に格納されている。この詳細は後で説明する。ステップ803では、図7で説明したストレージネットワークテーブルを送信先識別子をキーとしてサーチする。

【0039】ステップ804ではステップ803でサーチした結果を参照し、送信先がホストであればステップ805に移り、送信先がディスクアレイ制御装置であればステップ806に移る。ステップ805では、ホスト内のネットワーク処理プログラム405と同様の処理を行う。ステップ806では、専用高速送信処理プログラムの呼び出しを行う。この処理は後で詳細フローを示す。ステップ807は受信処理時に実行され、データに付加された情報からプロトコルを取得する。本実施例では、ホストから転送された場合は0が、ディスクアレイから転送された場合には1が格納されている。データ構造は後で詳細に説明する。

【0040】ステップ808では、送信元がホストかどうか判断する。この判断では、ステップ807で取得したプロトコルから判断可能である。ホストから転送された場合にはステップ809に移り、ディスクアレイ制御装置から転送された場合にはステップ810に移る。ステップ809では、ホスト内のネットワーク処理プログラム405と同様の処理を行う。ステップ810では、専用高速受信処理プログラムの呼び出しを行う。この処理は後で詳細フローを示す。このプロトコル選択プログラムにより、汎用インタフェースを要求する場合には、従来のネットワーク処理を行い、性能を要求する場合には専用高速送信/受信に切り替えることができる。

【0041】図9はネットワーク処理に使用されるデータ形式を示している。本データ構造は、ホストとディスクアレイ制御装置との間で通信を行う際に、ストレージネットワークを流れるデータの構造を示している。ホストからディスクアレイ制御装置に転送されたデータや、ディスクアレイ制御装置がホストに対してデータを送信する場合には図9のようなデータ形式で通信を行う。

【0042】901はストレージネットワークヘッダであり、ストレージネットワーク専用の制御情報が先頭に付加される。907はネットワークヘッダであり、ホスト内のネットワーク処理プログラム405が認識可能なデータのヘッダが付加される。912はユーザのコマンドあるいはデータである。ネットワークでは、図9に示すように階層を渡る毎にヘッダが付加される。送信側では、ネットワーク処理プログラムがまずネットワークヘ

ッダ907を付加し、ストレージネットワークを通過する際ストレージネットワークヘッダ901を付加する。受信側では、まずストレージネットワークヘッダ901が取り除かれ、次にネットワーク処理プログラムでネットワークヘッダ907が取り除かれ、ユーザプログラムに渡される。これはカプセル化と呼ばれており、各ヘッダを付加するために処理オーバーヘッドがかかるが、ユーザプログラムを変更する必要があることから一般的に採用されている方法である。

【0043】ネットワークヘッダ907内の908はユーザデータ長、909は送り元ネットワーク識別子、910は送り元ネットワーク識別子、911はエラーチェックコードが格納されている。908から910はデータを処理する上で必要な情報であるが、エラーチェックコード911はデータが正しく転送されたかどうかをチェックするために送信側でコードが計算され格納される。受信側では、送信側と同じようにエラーチェックコードを生成し転送されたエラーチェックコードと比較する。その結果、値が等しければデータは正しく転送されたと判断し、値が異なっていればデータの再送を送信側に要求する。

【0044】エラーチェックコードにはいくつかの種類があるが、一般的に用いられるのはチェックサムと呼ばれる方法で、送信するデータの内容を数バイト単位に区切り、和をとる方法である。この方法は、エラーチェック方法の中でも処理が簡単ではあるが、エラーチェックコードを生成するのにユーザデータ長に比例した時間がかかるため大きなオーバーヘッドとなる。

【0045】上記と同様にストレージネットワークヘッダ901内の902はプロトコル識別子、903はネットワークヘッダ907とユーザデータ&コマンド912の長さの合計、904は送り元識別子、905は送り先識別子、906はエラーチェックコードが格納されている。プロトコル識別子902は、ストレージネットワークヘッダ901の後に続くデータがどのようなプロトコルで転送されたかを示す識別子で、ネットワーク処理に使用される場合には0が格納される。エラーチェックコード906は、ネットワークヘッダ907で述べたエラーチェックコード911と同様である。機能としてはネットワークヘッダ907のエラーチェックと同様であり、2重のエラーチェックで信頼性を向上させる効果はあるが、処理の上ではオーバーヘッドが大きくなり転送性能は低下する。従来のネットワークもストレージネットワークも共に規格化されたプロトコルであり、またユーザインタフェースを変えない場合にこのようなオーバーヘッドが生じる。

【0046】図10は専用高速送信/受信の場合のデータ構造を示している。1001はストレージネットワークヘッダであり、1007はユーザデータ&コマンドである。ストレージネットワークヘッダ1001内の10

10

20

30

40

50

02はプロトコル識別子、1003はユーザデータ&コマンド1007の長さの合計、1004は送り元識別子、1005は送り先識別子、1006はエラーチェックコードが格納されている。

【0047】プロトコル識別子1002は、ストレージネットワークヘッダ1001の後に続くデータがどのようなプロトコルで転送されたかを示す識別子で、専用高速送信/受信の場合は1が格納される。これは受信したデータが、ホストから転送されたのかディスクアレイ制御装置から転送されたのかを判断するのに用いられる。専用高速送信/受信の場合は、図9の場合と異なりネットワークヘッダ907がない。ネットワーク形式のデータ構造ではないため、この形式ではホストと通信することはできないが、エラーチェックコードを2重に生成する必要がないため処理を高速化できる。

【0048】図11は、キャッシュ112、113からデータバッファ505に転送されたデータの構造を示している。1101はキャッシュヘッダ、1104はユーザデータ&コマンドが格納されている。このデータはマイクロプログラム118、119により生成される。キャッシュヘッダ1101内の1102は送り先識別子、1103はオプションにより付加されることもされないこともあるが、通常データ長等が格納されている。送り先識別子1102を参照することで、プロトコル選択プログラム(図8)内のステップ803で、どこにデータやコマンドを送信するかを判断することができる。

【0049】図12は専用高速受信プログラムの処理フローを示している。ステップ1201ではエラーチェックコードの確認を行う。これは、受信したデータのエラーチェックコードを生成し、図10のエラーチェックコード1006と比較することで実現可能である。ステップ1202ではエラーが発生したかどうかを判断する。エラーが発生していた場合、ステップ1207に移り、送信元に対し再送を要求し処理を終了する。エラーがなかった場合、ステップ1203に移る。ステップ1203では、送り元識別子を取り出す。ステップ1204では、ストレージネットワークヘッダを削除し、代わりにネットワークヘッダを付加する。ステップ1205では、キャッシュヘッダ内の送り先識別子1102にステップ1203で取り出した送り先識別子を格納する。これは処理終了後、ユーザデータやコマンドを返送する際に使用される。ステップ1206では、ステップ1205で生成したデータをキャッシュへ転送する。ディスクアレイ制御装置は、キャッシュにデータが転送されると転送された内容を解釈し実行する。

【0050】このように専用高速受信は、ネットワーク処理と異なり、エラーチェックコードチェック回数を削減することができるため処理を高速化できる。

【0051】図13は専用高速送信処理プログラムの処理フローを示している。ステップ1301では、キャッ

シュヘッダから送り先識別子の取り出しを行う。ステップ1302では、ストレージネットワークヘッダ内のプロトコル識別子を1に設定する。これは、専用高速送信データであることを受信側に知らせる意味を持つ。ステップ1303では、ストレージネットワークヘッダにデータ長を格納する。ステップ1304では、ストレージネットワークヘッダに送り先識別子を格納する。これはステップ1301でキャッシュヘッダから取得した識別子である。ステップ1305では、送り元識別子を格納する。ステップ1306ではエラーチェックコードを生成し、ストレージネットワークヘッダに格納する。このように専用高速送信は、ネットワーク処理と異なり、エラーチェックコードチェック回数を削減することができるため処理を高速化できる。

【0052】上記処理により、ディスクアレイ制御装置間でデータを高速に送受信することができるようになる。また、ユーザはこれらの切り替えは意識する必要がない。

【0053】一般にネットワークは、バスあるいはループのトポロジの意味が強いが、本発明はトポロジには依存しない。そのため、スイッチを使用したスター形式のトポロジでも発明の効果は得られる。

【0054】

【発明の効果】本発明では、複数のディスクアレイ制御装置をネットワークで接続した形態のストレージシステムにおいて、ホストとディスクアレイ制御装置間における通信プロトコルと、ディスクアレイ制御装置間における通信プロトコルを切り替える機構を設けることにより、ホストとは汎用インタフェースを保ちつつ、ディスクアレイ制御装置間では高速に通信を行うことができる。

【図面の簡単な説明】

【図1】本発明における実施例の全体のブロック図。

【図2】本発明の実施例におけるネットワークストレージのブロック図。

【図3】本発明の実施例における動作概要を示すブロック図。

【図4】本発明の実施例におけるホストプロセッサのブロック図。

【図5】本発明の実施例におけるストレージネットワーク制御装置のブロック図。

【図6】本発明の実施例におけるストレージネットワークのブロック図。

【図7】本発明の実施例におけるストレージネットワークテーブルを示す説明図。

【図8】本発明の実施例におけるプロトコル選択プログラムの処理のフローチャート。

【図9】本発明の実施例におけるネットワークプロトコル転送時のデータ構造を示す説明図。

【図10】本発明の実施例における専用高速送信/受信プロトコル転送時のデータ構造を示す説明図。

13

【図11】本発明の実施例におけるキャッシュヘッダの構造を示す説明図。

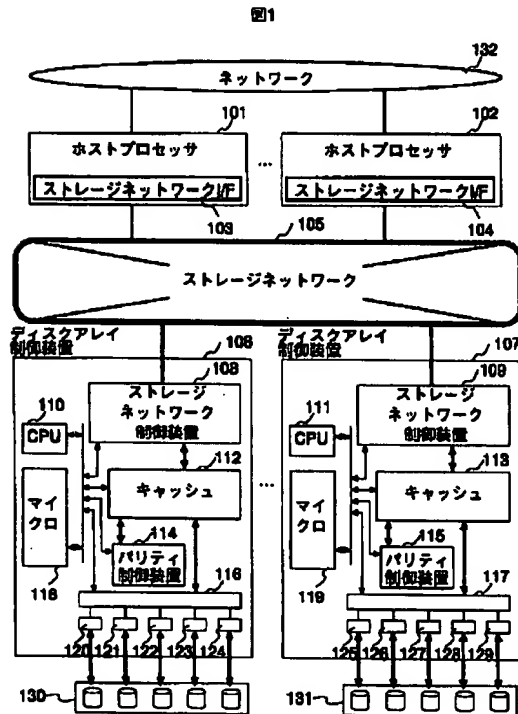
【図12】本発明の実施例における専用高速受信プログラムのフローチャート。

【図13】本発明の実施例における専用高速送信プログラムのフローチャート。

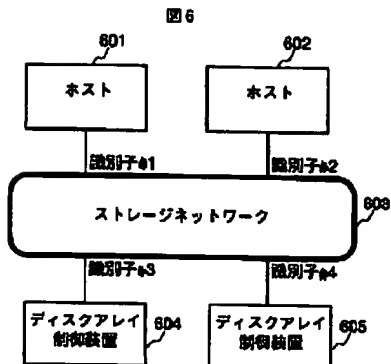
【符号の説明】

*

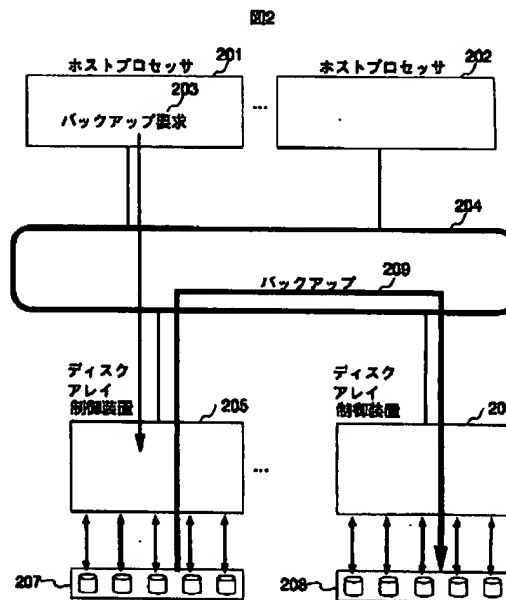
【図1】



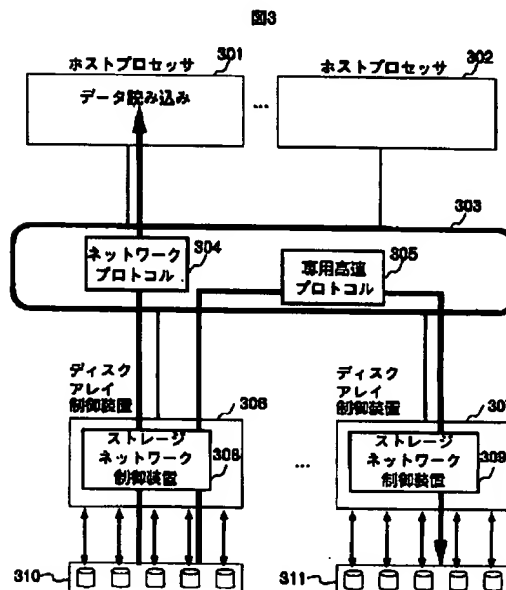
【図6】



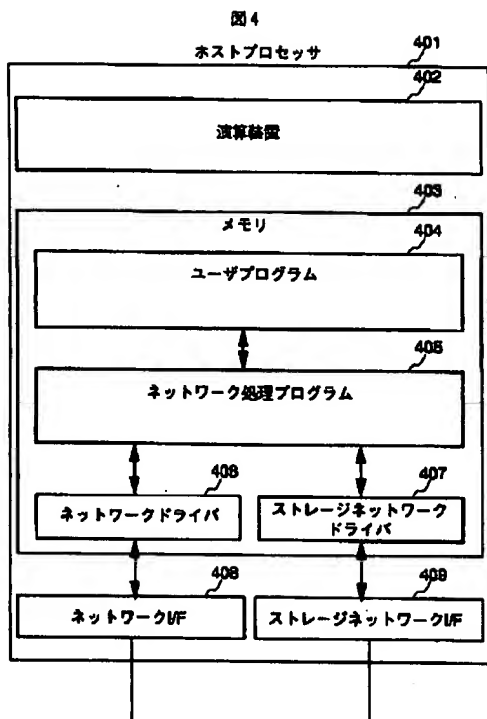
【図2】



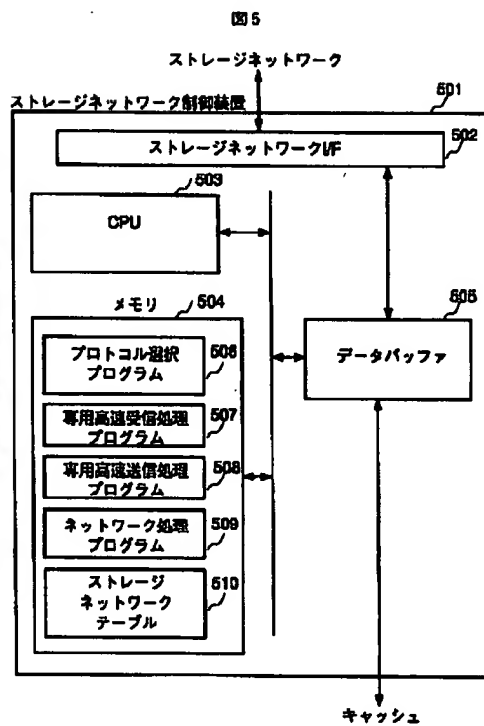
【図3】



【図4】



【図5】

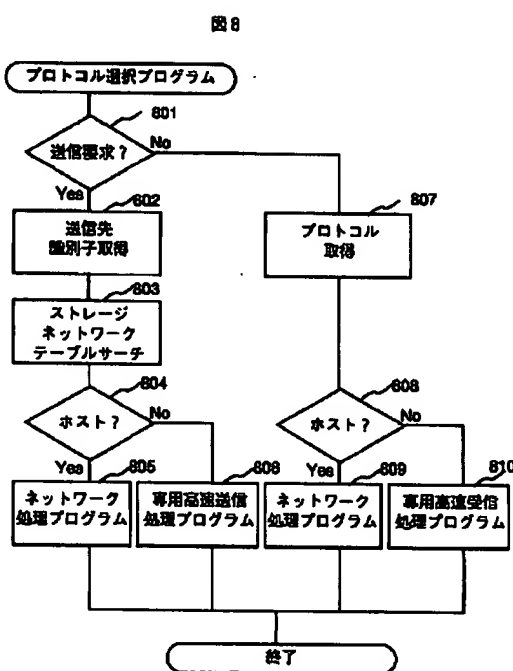


【図7】

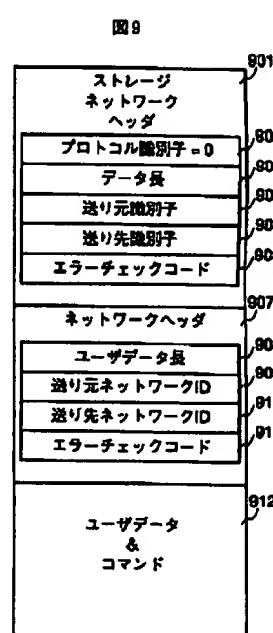
図7

ストレージネットワーク識別子 701	装置種別 702
1	0 (ホスト)
2	0 (ホスト)
3	1 (制御装置)
4	1 (制御装置)

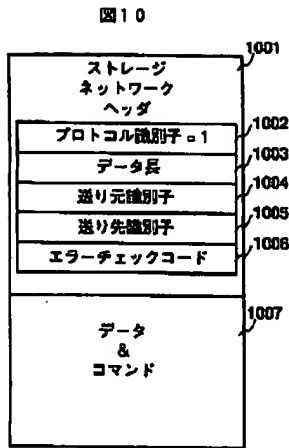
【図8】



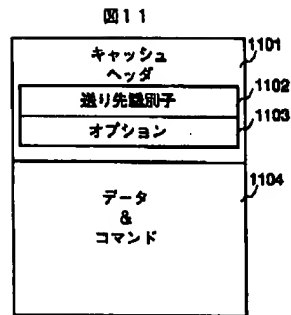
【図9】



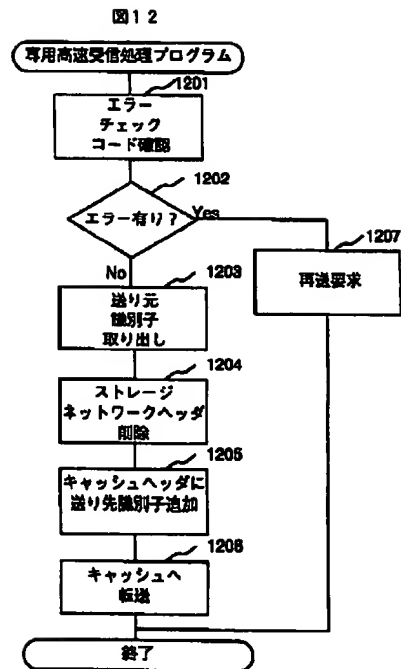
【図10】



【図11】



【図12】



【図13】

